**Big Data and Predictive Analytics
in Business Use Cases**

HAMBURG,
March 10, 2016

**SIM** + **COG**

SIMULATION · COGNITION

**SimCog Technologies GmbH**

- Founded in 2012

- Based in Hamburg

- 6 Employees → Looking for more Data Scientists

- Several Use Cases in the Area of Predictive Analytics

- Mainly Private Equity owned

SIMCOG
TECHNOLOGIES

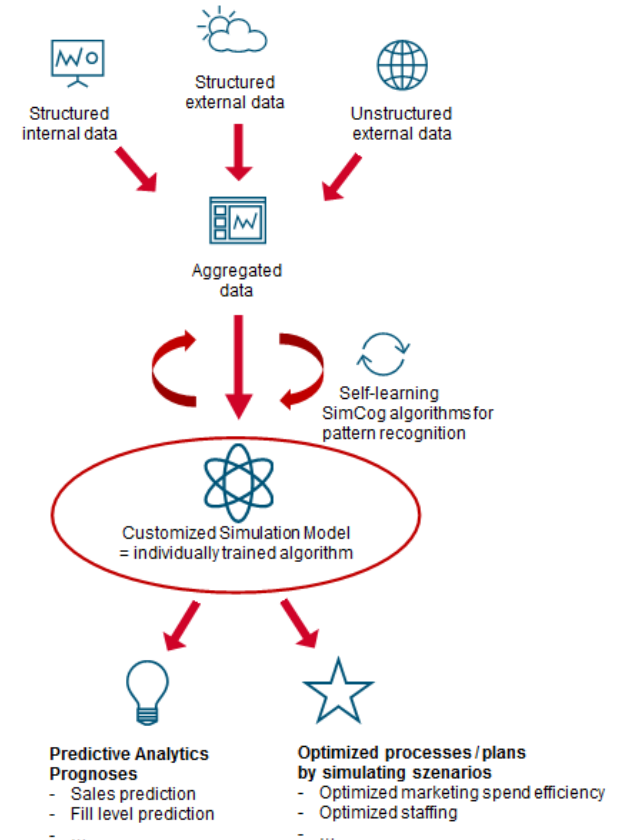**Difference between Science and Business**

- Shorter turnaround-time in business

    - Projects last ~ 3 Month

- Specialized "skill set" of Data Scientists

    - Finding economic relevant problem and external data (Consultant)

    - Explaining methods / results / tests to "uneducated" managers  (Sales / Acquisition)

    - Find best statistical method, e.g. Neural Net vs. BDT (Statistician)

    - Efficient programming in small teams (Programmer)

    - Extract relevant data from Databases, APIs, Web crawler, … (Data Analyst)

- Long-Term Jobs and better Payment

SIMCOG
TECHNOLOGIES

## Fusion of relevant Data

- Structured internal company data
- Structured external data, e.g. geodata, holiday seasons, weather, etc.
- Special feature: Inclusion of unstructured data is possible, e.g. social media monitoring

## Properties of the Simulation Model

- Self-learning algorithms
- Quality checks with historic data
- Selection of appropriate algorithms for specific questions



Structured internal data

Structured external data

Unstructured external data

Aggregated data

Self-learning SimCog algorithms for pattern recognition

Customized Simulation Model = individually trained algorithm

**Predictive Analytics Prognoses**
- Sales prediction
- Fill level prediction
- …

**Optimized processes / plans by simulating szenarios**
- Optimized marketing spend efficiency
- Optimized staffing
- …

## What does "Big Data" mean in Research and Business

| | | | |
|---|---|---|---|
| Byte | = | Grain of Rice | |
| Kilobyte | = | Cup of Rice | |
| Megabyte | = | Eight Bags of Rice | |
| **Gigabyte** | **=** | **Three Lorries of Rice** | → **Business** |
| Terabyte | = | Two Container ships of Rice | |
| **Petabyte** | **=** | **Manhattan covered with Rice** | → **CERN** |
| Exabyte | = | Great Britain covered with Rice (3 times) | |
| Zettabyte | = | Fills Pacific Ocean with Rice | |

**However, data in business comes from many different sources and the inclusion of external data is the key in big data analysis**

SimCog
TECHNOLOGIES

## The Good, the Bad and the Ugly

**Some Data-Format Examples**

- Most Data available in „good" CSV-files
  - Inconsistent Data/Time Formats
  - , and . in numbers
  - Missing quotation and binary numers

- Database Dump

```
KUNDEN_NR_        POSITION_NR_     STATUS  Kundenart        Kundengruppe    ADM      Partnertyp        PREISGRUPPE
MINERALOLSTEUER_KZ               Branche ANWENDUNGSART    VORNAME STRASSE LAND_KZ_   POSTLEITZAHL   ORT      VERTRAGART
VERTRAGDATUM     Miete   Miete Zeiteinheit        Wartung Wartung Zeiteinheit       Abschlag       Abschlag Zeiteinheit
LIEFERDATUM      ARTIKEL_NR_      BEZEICHUNG_1     MENGE_KG          MENGE_LITER      MENGENEINHEIT  PROZ_NACH_BEFUELLUNG
TEMPERATUR       POSITIONSWERT_NETTO      BEZEICHNUNG      BEZEICHNUNG2     Beh.-Gr"áe to   Lagerungsart      LFD_BEHAELTER_NR
beh__nr__hersteller       hersteller       baujahr behaltervolumen
77005354         1         I       Z„hleranlage             Z„hlerkunden             V                NULL    63      BT
unbekannt                 Heizen und Kochen                                                                           D
                          24plus-Miete                      01.01.1960
02.10.2001       110015  BRENNGAS-Tank                                              1539,462         3005   LTR     60      15
998,68  100 Liter                 Sofort netto f„llig       2,1 To            unterirdisch           2714100 82287  18
01.01.1998       4850
77005349         1         I       Tank-Endverbraucher      Tank Endverbraucher                               Betreiber
         0       BT      unbekannt        Heizen und Kochen        VW                                61,36  H
         D                                Tank-Miete/Wartung       01.01.1960                                2769
47,55   H                         25.09.2003       110015  BRENNGAS-Tank                                      oberirdisch
5364    LTR     84      10      2708,82 100 Liter                 Sofort netto f„llig       2,9 To
```

SimCog
TECHNOLOGIES

## The Good, the Bad and the Ugly

**Some Data-Format Examples**

- Most Data available in „good" CSV-files
    - Inconsistent Data/Time Formats
    - , and . in numbers
    - Missing quotation and binary numers

- Database Dump

- XML in CSV-file

| 3 | <NM_DOCS xmlns="http://api.nexmart.net/services/nexmartmeta/2014-05-25" x |
| 4 | <NM_DOC> |
| 5 | <DOCUMENT type="orders" qualifier="default" role="original" test="false"> |
| 6 | <VERSION>4.0</VERSION> |
| 7 | <HEADER> |
| 8 | <CONTROL_INFO> |
| 9 | <LAST_SAVE_DATE>2015-10-26T08:38:39+01:00</LAST_SAVE_DATE> |
| 10 | <PROCESS_TYPE>silent</PROCESS_TYPE> |
| 11 | <SOURCE>edi</SOURCE> |
| 12 | <DESTINATION████████</DESTINATION> |
| 13 | <LOGS> |
| 14 | <LOG type="buyer"> |
| 15 | <MESSAGE_ID>FE020/OE002/0247</MESSAGE_ID> |
| 16 | </LOG> |
| 17 | <LOG type="nm_in"> |
| 18 | <MESSAGE_ID>75704088</MESSAGE_ID> |
| 19 | <PROCESS_DATE>2015-10-26T08:38:31+01:00</PROCESS_DATE> |
| 20 | </LOG> |
| 21 | <LOG type="nm_messageprocess"> |
| 22 | <MESSAGE_ID>OP-0001768f-67c8-4fe7-8dac-55557b7e0e8c</MESSAGE_ID> |
| 23 | <PROCESS_DATE>2015-10-26T08:38:31+01:00</PROCESS_DATE> |
| 24 | </LOG> |
| 25 | </LOGS> |

SimCog
TECHNOLOGIES

## The Good, the Bad and the Ugly

**Some Data-Format Examples**

- Most Data available in „good" CSV-files
    - Inconsistent Data/Time Formats
    - , and . in numbers
    - Missing quotation and binary numers

- Database Dump

- XML in CSV-file

- Excel Sheets



| | A | B | C | D | E | F | G | H | I | J | K | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2 | Gesamt Online ( Franchise und eigene Standorte ) | | | | | | | | | | | |
| 3 | Home- und Landingpages (Google/Bing) | | | | | | | | | | | |
| 4 | | | | | | | 2015 | | | | | |
| 5 | | | | | | | | | | | | |
| 6 | Datum | Tag | Budget SEA ( Google + Bing ) | Visits | Mails | Mails über | VCR Mail | Anrufe | Anrufe über | VCR Anrufe | Anfragen gesamt | VCR Gesamt |
| 216 | 29.07.2015 | Mittwoch | 2,509 € | 3,414 | 65 | 0 | 1.9% | 95 | 0 | 2.8% | 160 | 4.7% |
| 217 | 30.07.2015 | Donnerstag | 2,061 € | 3,051 | 51 | 0 | 1.7% | 93 | 0 | 3.0% | 144 | 4.7% |
| 218 | 31.07.2015 | Freitag | 1,675 € | 2,397 | 45 | 0 | 1.9% | 78 | 0 | 3.3% | 123 | 5.1% |
| 219 | 01.08.2015 | Samstag | 787 € | 1,217 | 26 | 0 | 2.1% | 10 | 0 | 0.8% | 36 | 3.0% |
| 220 | 02.08.2015 | Sonntag | 967 € | 1,479 | 22 | 0 | 1.5% | 14 | 0 | 0.9% | 36 | 2.4% |
| 221 | 03.08.2015 | Montag | 1,772 € | 2,843 | 47 | 0 | 1.7% | 77 | 0 | 2.7% | 124 | 4.4% |
| 222 | 04.08.2015 | Dienstag | 1,809 € | 2,683 | 39 | 0 | 1.5% | 69 | 0 | 2.6% | 108 | 4.0% |
| 223 | 05.08.2015 | Mittwoch | 1,459 € | 2,347 | 33 | 0 | 1.4% | 54 | 0 | 2.3% | 87 | 3.7% |
| 224 | 06.08.2015 | Donnerstag | 1,195 € | 2,131 | 31 | 0 | 1.5% | 57 | 0 | 2.7% | 88 | 4.1% |
| 225 | 07.08.2015 | Freitag | 1,121 € | 1,810 | 18 | 0 | 1.0% | 33 | 0 | 1.8% | 51 | 2.8% |
| 226 | 08.08.2015 | Samstag | 523 € | 975 | 14 | 0 | 1.4% | 11 | 0 | 1.1% | 25 | 2.6% |
| 227 | 09.08.2015 | Sonntag | 1,003 € | 1,384 | 22 | 0 | 1.6% | 6 | 0 | 0.4% | 28 | 2.0% |
| 228 | 10.08.2015 | Montag | 1,769 € | 2,759 | 45 | 0 | 1.6% | 62 | 0 | 2.2% | 107 | 3.9% |
| 229 | 11.08.2015 | Dienstag | 1,947 € | 2,920 | 34 | 0 | 1.2% | 55 | 0 | 1.9% | 89 | 3.0% |
| 230 | 12.08.2015 | Mittwoch | 1,939 € | 3,274 | 21 | 0 | 0.6% | 78 | 0 | 2.4% | 99 | 3.0% |
| 231 | 13.08.2015 | Donnerstag | 1,824 € | 2,699 | 41 | 0 | 1.5% | 30 | 0 | 1.1% | 71 | 2.6% |
| 371 | | | | | | | | | | | | |
| 372 | | | | | | | | | | | | |
| 373 | | | | | | | | | | | | |
| 374 | | | | | | | | | | | | |
| 375 | Die VCR ist erst ab dem 09.01.2015 korrekt mit dem Vorjahr vergleichbar | | | | | | | | | | | |
| 376 | Die Visits wurden auf Basis der Entwicklung (+13% Visits) in 2015 für 2014 korrigiert. | | | | | | | | | | | |

# The Good, the Bad and the Ugly

## Some Data-Format Examples

- Most Data available in „good" CSV-files
  - Inconsistent Data/Time Formats
  - , and . In numbers
  - Missing quotation and binary numers

- Database Dump

- XML in CSV-file
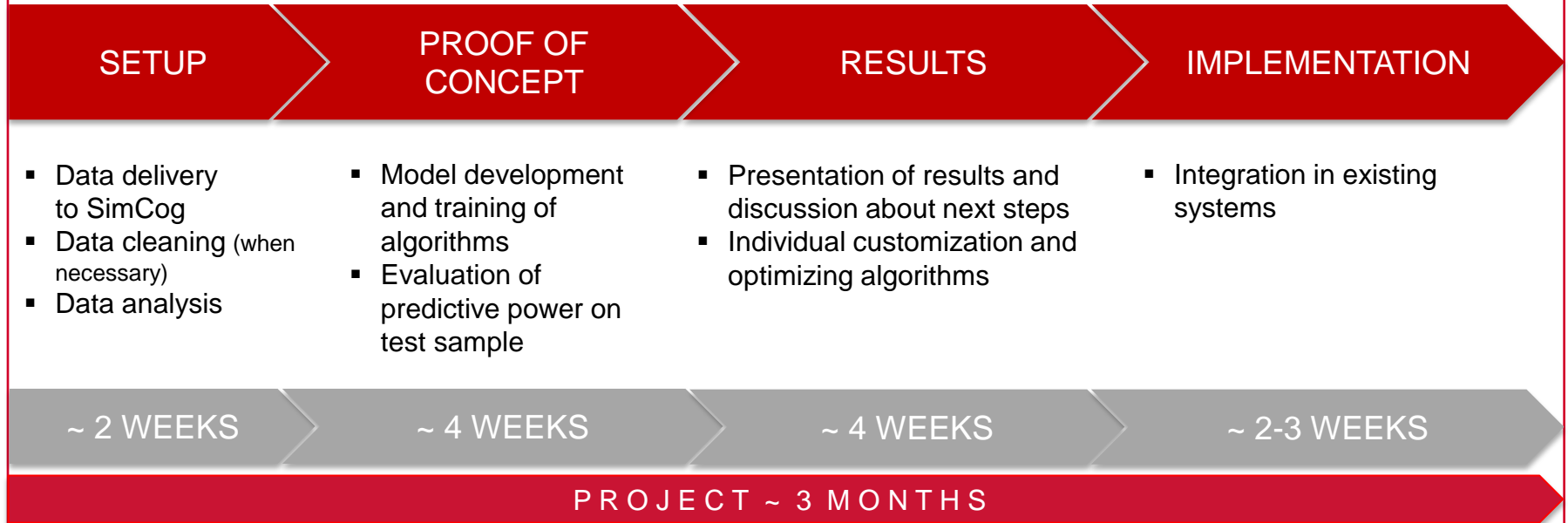
- Excel Sheets

- „Media Plan"-Excel Sheets

- PDF

**Programming Languages**

- C++

- ROOT    (only at SimCog)

- Java (Natural Language Processing, APIs)

- R

- Python

- Database (SQL, noSQL etc.)

SIMCOG
TECHNOLOGIES

**First Steps before Starting a Project**

- Task: Find an economic relevant problem, that can be answered with the company data*
- Required prediction quality depends on economic leverage

| SETUP | PROOF OF CONCEPT | RESULTS | IMPLEMENTATION |
|---|---|---|---|
| - Data delivery to SimCog<br>- Data cleaning (when necessary)<br>- Data analysis | - Model development and training of algorithms<br>- Evaluation of predictive power on test sample | - Presentation of results and discussion about next steps<br>- Individual customization and optimizing algorithms | - Integration in existing systems |
| ~ 2 WEEKS | ~ 4 WEEKS | ~ 4 WEEKS | ~ 2-3 WEEKS |

P R O J E C T ~ 3 M O N T H S

*Typically companies have historical data of the past ~ 4 years

SIMCOG
TECHNOLOGIES

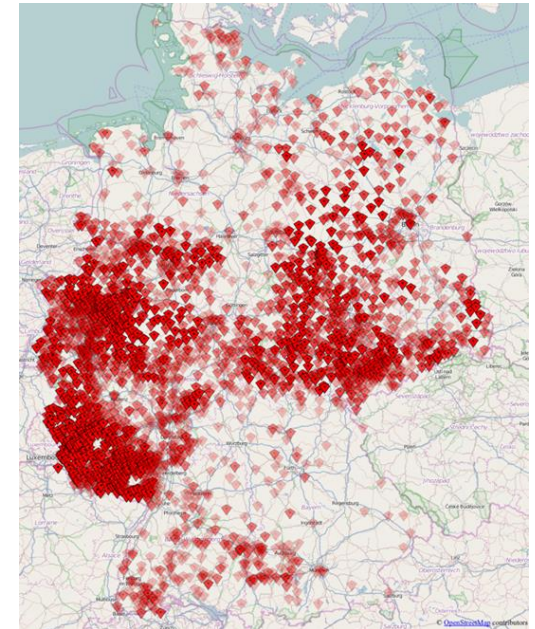USE CASE
ENERGY

**Initial question**

- Client is an energy supplier that delivers gas to ~ 30.000 consumption points in Germany
- Basic idea: optimize delivery logistics by a precise fill level prediction
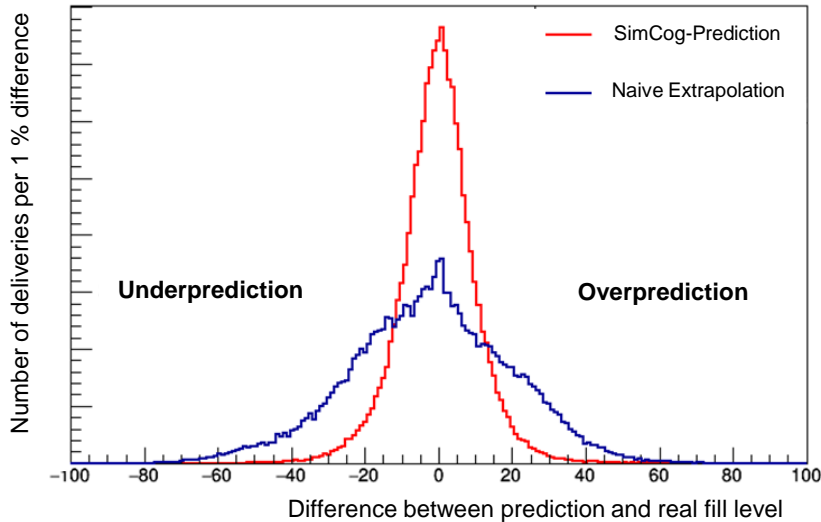
**Used internal data** → *provided by client*

- All gas deliveries since 2002
  - Quantity delivered, date of delivery, fill level after filling, …
- Contract type
  - Heating, heating & cooking, balloonist, industry, …
- Fuel tank capacity
- Use of alternative energy sources
- Postal code

**Used external data** → *added by SimCog*

- Weather (per consumption point weather information from the three nearest weather stations is used)
- Geo-structural data, e.g. population density
- Information of the tank supplier

SimCog
TECHNOLOGIES

USE CASE
ENERGY

Comparison Naive Extrapolation / SimCog-Prediction



**Comparison of predictions:**

The method „Naive Extrapolation" considers the past and carries on the average consumption of the past

The SimCog-Prediction uses latest machine-learning-algorithms and integrates external data

→ SimCog´s pattern recognition is more than twice as accurate as the naive extrapolation

→ **Precision**
Very precise predictions: uncertainty ~ 5% on sales volumes

SIMCOG
TECHNOLOGIES

USE CASE
ENERGY

- Fill level prediction
  - Daily updates of the fill level prediction for 30'000 existing customers

- Delivery date prediction
  - Daily updates of the delivery date prediction for 30'000 existing customers

- Identification of illegal third-party refillment

- Quarterly prediction for optimized supply chain management, energy disposition, purchasing & liquidity planning (of the client)

- Detection of essential data errors
  - Issue warnings for unexpected events such as: negative consumption, illegal third-party refillment, strong behavior change, data error, …

- Automatic adjustment of predictions when behavior is altered (self-learning)

SimCog
TECHNOLOGIES

USE CASE
ENERGY

1. Reduced storage and standby costs

   - Optimized supply chain management through greatly improved estimated amount of energy needed

2. Cost savings through improved route planning

   - Optimized distribution logistics

   - Better utilization of refilling trucks

   - Customers can be addressed directly for e.g. „early refueling"

3. Improved marketing and sales activities

   - Improved timing of customer contact

   - Special promotion more controllable

4. Improved handling of illegal third-party refillment

   - Phone customers where the probablility of third-party refillment is high

   - Legal action is an option in safe cases

5. Improved overview of turnover by predicting the fill level in the counter systems

SIMCOG
TECHNOLOGIES

## Retail

### Retail B2C

- Sales Prediction
- Footfall
- Customer Forecast
- Shopping Cart
- Online / Offline

### Retail B2B

- Sales Prediction
- Online / Offline

## Marketing

### Marketing

- Marketing Spend Efficiency
- Retargeting
- Churn Rate
- Coupon-Conversion

## Logistics

### Logistics

- Shipping ETA
- Fill Level Prediction Energy

## More Solutions

### More Solutions

- Fraud Detection

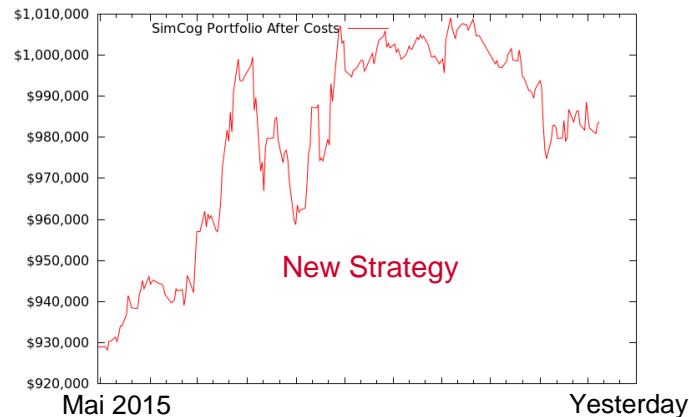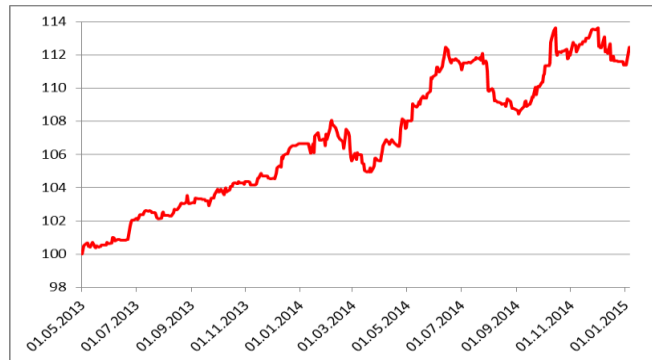- Stock Price Prediction

- Protection against Economic Espionage

- …

SimCog
TECHNOLOGIES

USE CASE
STOCK PRICE

## Predictive Analytics: Stock Price

Result of real trades on SimCog account:



**Use Case**

- Predict most likely stock price movements and trade with an automated market-neutral trading strategy

- In addition to using stock prices, sector information, director dealings, etc. the analysis of social media information is a crucial factor
    - Discussion from the social web are considered within individual subject areas

SimCog
TECHNOLOGIES

**Dr. Jan Thomsen**
Managing Director

**SimCog Technologies GmbH**
Schauenburgerstr. 27
20095 Hamburg

T: +49 40 2000 396 07
F: +49 40 2000 396 20
E: thomsen@simcog.de

www.simcog.de

SIMCOG
TECHNOLOGIES